



INFERENCIA Y MODELOS ESTADÍSTICOS

Jacqueline Köhler C. y José Luis Jara V.



CAPÍTULO 6. PODER ESTADÍSTICO

En el capítulo 4 estudiamos el procedimiento para someter hipótesis a prueba, junto con los errores de decisión que podríamos cometer:

- Error tipo I: rechazar H_0 en favor de H_A cuando H_0 es en realidad verdadera.
- Error tipo II: no rechazar H_0 en favor de H_A cuando H_A es en realidad verdadera.

Allí conocimos el nivel de significación, α , como herramienta para representar y, de alguna manera, controlar la probabilidad de cometer un error de tipo I, con lo que la preocupación se centra en controlar la ocurrencia de esta clase de errores, desviando la atención de los errores de tipo II. Esto se debe a que la hipótesis nula representa el *status quo*, es decir, mantener las cosas y creencias tal como están y, por ende, cuando no se rechaza H_0 , no suele requerirse tomar ninguna acción. En contraste, la hipótesis alternativa describe un cambio de condiciones, por lo que rechazar H_0 en favor de H_A usualmente conlleva un esfuerzo, mayor costo, para adaptarse o aprovechar las nuevas condiciones.

Sin embargo, en el capítulo 4 también vimos que el valor de α debe ser acorde con las consecuencias de cometer errores tanto de tipo I como de tipo II, ¡pero no sabemos cómo se relaciona el nivel de significación con los errores de tipo II!

Así como el nivel de significación α corresponde a la probabilidad de cometer errores de tipo I, definimos ahora β como la probabilidad de cometer errores de tipo II. α y β están relacionados: **para un tamaño fijo de la muestra: al reducir β , α aumenta, y viceversa**. Este fenómeno se evidencia con mayor fuerza mientras más pequeña sea la muestra. No obstante, en la práctica resulta más interesante conocer la probabilidad de **no** cometer errores de tipo II. Esto nos lleva a un nuevo concepto: el **poder estadístico** de una prueba de hipótesis, dado por $1 - \beta$, que se define como **la probabilidad de correctamente rechazar H_0 cuando es falsa**.

Otra forma de entender la noción de poder de una prueba es qué tan propensa es esta para distinguir un efecto real de una simple casualidad, lo que nos lleva a la noción de **tamaño del efecto**, que corresponde a una cuantificación de la diferencia entre dos grupos, o del valor observado con respecto al valor nulo.

En el capítulo 5 conocimos la prueba t para inferir acerca de dos medias. En este contexto, el tamaño del efecto corresponde a qué tan grande es la diferencia real entre ambas. Si quieres aprender más sobre estos conceptos, puedes consultar las fuentes en las que se basa este capítulo: Diez y col. (2017, pp. 239-245) y Freund y Wilson (2003, pp. 123-138).

6.1 PODER, NIVEL DE SIGNIFICACIÓN Y TAMAÑO DE LA MUESTRA

En la introducción de este capítulo vimos que el poder corresponde a la probabilidad de **no** cometer un error de tipo II, y que está muy relacionado con el tamaño de la muestra. También mencionamos que existe una relación entre el poder y el nivel de significación, la cual exploraremos en esta sección.

La figura 6.1 (creada mediante el script 6.1) muestra cuatro curvas de poder para la prueba t de Student de una muestra con desviación estándar $s = 1$ y valor nulo $\mu_0 = 0$. En ella, el tamaño del efecto está representada en la misma escala de la variable, aunque en la sección siguiente veremos otra alternativa. La curva roja considera $\alpha = 0,05$ y $n = 6$; la azul, $\alpha = 0,01$ y $n = 6$; la verde, $\alpha = 0,05$ y $n = 10$, y la naranja, $\alpha = 0,01$ y $n = 10$. En ella podemos observar que:

- El poder de la prueba aumenta mientras mayor es el tamaño del efecto (en este caso, la distancia entre el valor nulo y la media de la muestra).
- A medida que el tamaño del efecto disminuye (es decir, el estimador se acerca al valor nulo), el poder se aproxima al nivel de significación.
- Usar un valor de α más exigente (menor), manteniendo constante el tamaño de la muestra, hace que la curva de poder sea más baja para cualquier tamaño del efecto (lo que verifica la relación entre α y β).
- Usar una muestra más grande aumenta el poder de la prueba para cualquier tamaño del efecto distinto de 0.

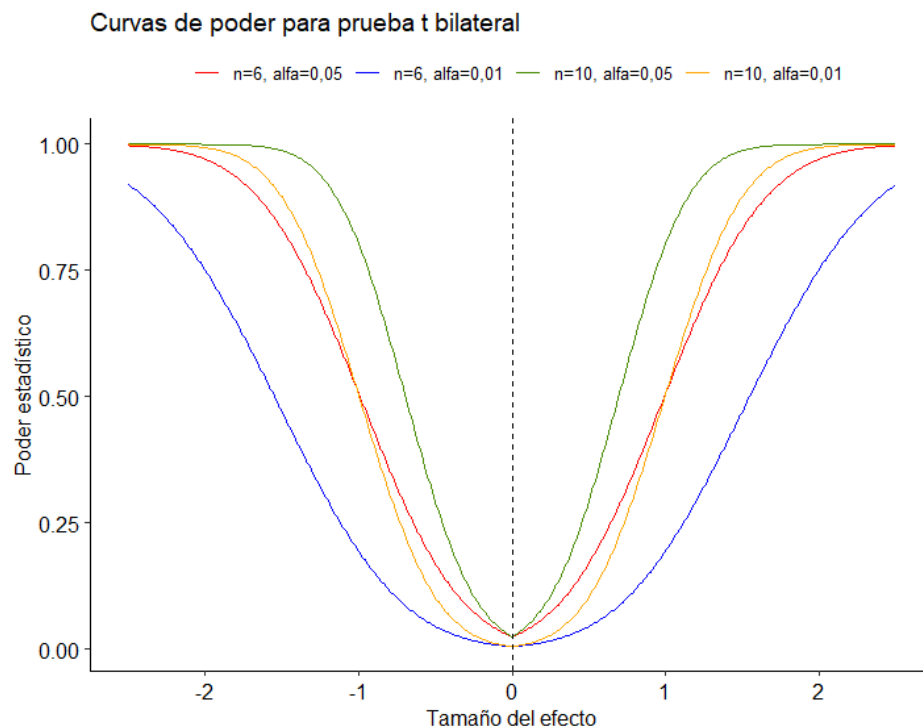


Figura 6.1: poder estadístico para prueba t bilateral.

De manera similar, la figura 6.2 considera las mismas muestras y los mismos niveles de significación que la figura 6.1, pero ahora para una prueba t unilateral. En ella se aprecia que la gran desventaja de las pruebas unilaterales es que el poder tiende a cero a medida que el tamaño del efecto aumenta en sentido contrario a la hipótesis alternativa, por lo que no sería posible detectar una diferencia en el sentido opuesto aunque fuese muy grande (pues no hay una región de rechazo en dicho sentido). El script empleado para la construcción de la figura 6.2 es idéntico al script 6.1, excepto porque el argumento `alternative` toma como valor “one.sided” en las llamadas a `power.t.test()`.

Script 6.1: poder estadístico para prueba t bilateral.

```
1 library(ggpubr)
2 library(tidyverse)
3
4 # Generar un vector con un rango de valores para la efecto
5 # de medias.
6 efecto <- seq(-2.5, 2.5, 0.01)
7
8 # Calcular el poder para una prueba t bilareral, para cada tamaño
9 # del efecto, asumiendo una muestra con desviación estándar igual a 1.
10 # Se consideran 4 escenarios para calcular el poder:
11 # 1. Una muestra de tamaño 6 y nivel de significación 0.05.
```

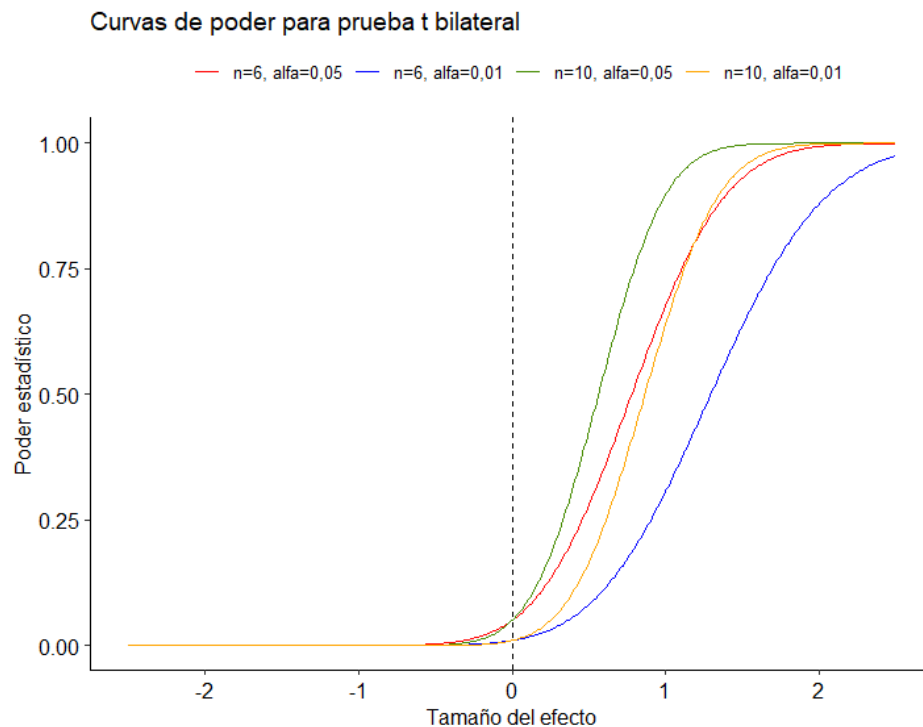


Figura 6.2: poder estadístico para prueba t unilateral.

```

12 # 2. Una muestra de tamaño 6 y nivel de significación 0.01.
13 # 3. Una muestra de tamaño 10 y nivel de significación 0.05.
14 # 4. Una muestra de tamaño 10 y nivel de significación 0.01.
15 n_6_alfa_05 <- power.t.test(n = 6,
16                             delta = efecto,
17                             sd = 1,
18                             sig.level = 0.05,
19                             type = "one.sample",
20                             alternative = "two.sided")$power
21
22 n_6_alfa_01 <- power.t.test(n = 6,
23                             delta = efecto,
24                             sd = 1,
25                             sig.level = 0.01,
26                             type = "one.sample",
27                             alternative = "two.sided")$power
28
29 n_10_alfa_05 <- power.t.test(n = 10,
30                             delta = efecto,
31                             sd = 1,
32                             sig.level = 0.05,
33                             type = "one.sample",
34                             alternative = "two.sided")$power
35
36 n_10_alfa_01 <- power.t.test(n = 10,
37                             delta = efecto,
38                             sd = 1,
39                             sig.level = 0.01,
40                             type = "one.sample",
41                             alternative = "two.sided")$power

```

```

42
43 # Construir matriz de datos en formato ancho.
44 datos <- data.frame(efecto, n_6_alfa_05, n_6_alfa_01,
45                       n_10_alfa_05, n_10_alfa_01)
46
47 # Llevar a formato largo.
48 datos <- datos %>% pivot_longer(!"efecto",
49                                names_to = "fuente",
50                                values_to = "poder")
51
52 # Formatear fuente como variable categórica.
53 niveles <- c("n_6_alfa_05", "n_6_alfa_01", "n_10_alfa_05",
54             "n_10_alfa_01")
55
56 etiquetas <- c("n=6, alfa=0,05", "n=6, alfa=0,01", "n=10, alfa=0,05",
57              "n=10, alfa=0,01")
58
59 datos[["fuente"]] <- factor(datos[["fuente"]], levels = niveles,
60                             labels = etiquetas)
61
62 # Graficar las curvas de poder.
63 g <- ggplot(datos, aes(efecto, poder, colour = factor(fuente)))
64 g <- g + geom_line()
65 g <- g + labs(colour = "")
66 g <- g + ylab("Poder estadístico")
67 g <- g + xlab("Tamaño del efecto")
68
69 g <- g + scale_color_manual(values=c("red", "blue", "chartreuse4",
70                                     "orange"))
71
72 g <- g + theme_pubr()
73 g <- g + ggtitle("Curvas de poder para prueba t bilateral")
74 g <- g + geom_vline(xintercept = 0, linetype = "dashed")
75
76 print(g)

```

La figura 6.3 muestra las curvas de poder para una prueba t unilateral y otra bilateral, ambas para una muestra de tamaño 6, desviación estándar $s = 1$ y $\alpha = 0,05$. En ella se evidencia claramente la ventaja de las pruebas unilaterales: cuando el tamaño del efecto aumenta en el sentido de la hipótesis alternativa, el poder es mayor que para una prueba bilateral.

Es deseable que las pruebas que se empleen para docimar hipótesis tengan un alto poder y, si hay más de una prueba disponible, se debe escoger la más poderosa. No obstante, los cálculos del poder suelen ser altamente complejos. Afortunadamente, la teoría permite en muchos casos conocer la prueba con mayor poder posible ante cualquier hipótesis alternativa, nivel de significación y tamaño de muestra (siempre que se cumplan las condiciones de base). Estas pruebas reciben el nombre de **uniformemente más poderosas**, y tal es el caso de la prueba t de Student.

6.2 TAMAÑO DEL EFECTO

El problema que podríamos tener al considerar el tamaño del efecto en la misma escala de la variable estudiada, como hemos hecho hasta ahora, es que esta escala varía de variable en variable. Para poder hacer comparaciones con mayor libertad, existen diferentes **medidas estandarizadas de efecto** que podemos

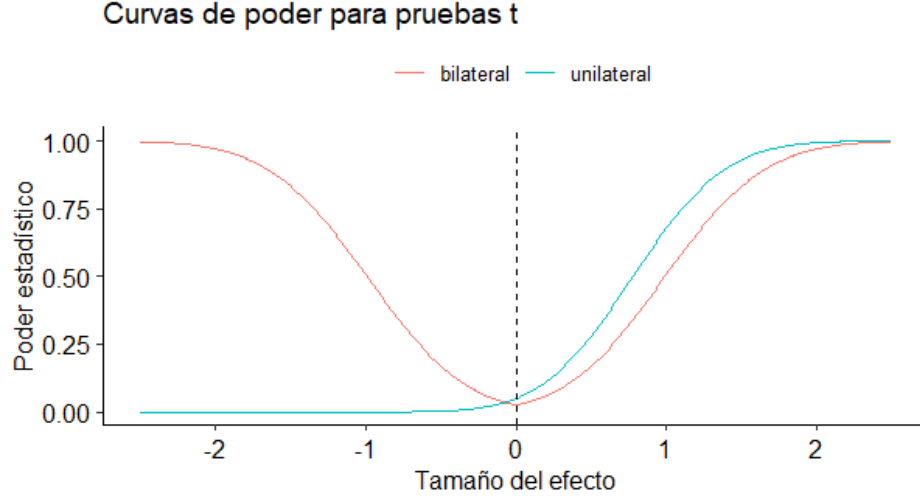


Figura 6.3: poder estadístico para pruebas t.

usar. Puesto que hasta ahora solo hemos estudiado la prueba t de Student, en esta sección conoceremos la llamada **d de Cohen** (Kassambara, 2019), una medida estándar ampliamente empleada para el tamaño del efecto con esta prueba.

En términos generales, se considera que $d = 0,2$ es un efecto pequeño (imperceptible a simple vista), $d = 0,5$ es un efecto mediano (probablemente perceptible a simple vista) y $d = 0,8$, un efecto grande (definitivamente perceptible a simple vista).

En el caso de la prueba t de una muestra, la d de Cohen se calcula como muestra la ecuación 6.1, donde:

- \bar{x} : media muestral.
- μ_0 : media teórica para el contraste (valor nulo).
- s : desviación estándar de la muestra con $n - 1$ grados de libertad.

$$d = \frac{\bar{x} - \mu_0}{s} \quad (6.1)$$

Para la prueba t de diferencia de dos medias (también llamada prueba t para dos muestras independientes o, simplemente, prueba t independiente), si el tamaño de la muestra es mayor a 50 elementos, se calcula como muestra la ecuación 6.2, y para muestras pequeñas se aplica un factor de corrección, como indica la ecuación 6.3, donde:

- \bar{x}_1, \bar{x}_2 : medias muestrales de cada grupo.
- n_1 y n_2 son los tamaños de ambas muestras.
- s_p : desviación estándar agrupada, dada por la ecuación 6.4¹.

$$d = \frac{\bar{x}_1 - \bar{x}_2}{s_p} \quad (6.2)$$

$$d = \frac{\bar{x}_1 - \bar{x}_2}{s_p} \cdot \frac{n_1 + n_2 - 3}{n_1 + n_2 - 2, 25} \quad (6.3)$$

$$s_p = \sqrt{\frac{\sum(x - \bar{x}_1)^2 + \sum(x - \bar{x}_2)^2}{n_1 + n_2 - 2}} \quad (6.4)$$

¹Note que esta corresponde a la raíz de la varianza agrupada descrita en 5.5

En el caso de la variante de Welch para la prueba t independiente, la fórmula para el cálculo de la d de Cohen es ligeramente diferente, como puede apreciarse en la ecuación 6.5.

$$d = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2 + s_2^2}{2}}} \quad (6.5)$$

Por último, las ecuaciones 6.6 y 6.7 muestran cómo se calcula la d de Cohen en el caso de la prueba t con muestras pareadas grandes ($n > 50$) y pequeñas, respectivamente, donde D corresponde a las diferencias entre las observaciones pareadas.

$$d = \frac{\bar{x}_D}{s_D} \quad (6.6)$$

$$d = \frac{\bar{x}_D}{s_D} \cdot \frac{n_1 - 2}{n_1 - 1, 25} \quad (6.7)$$

6.3 PODER, TAMAÑO DEL EFECTO Y TAMAÑO DE LA MUESTRA

Mencionamos en páginas anteriores que el poder puede también entenderse como qué tan propensa es una prueba estadística para distinguir un efecto real de una simple casualidad, y que podemos cuantificar este efecto.

Una gran ventaja del poder estadístico es que nos sirve para determinar el tamaño adecuado de la muestra para detectar un cierto tamaño del efecto. La figura 6.4, elaborada con el script 6.2, muestra el aumento del poder estadístico a medida que el tamaño de la muestra aumenta (para un tamaño del efecto y nivel de significación fijos). En ella se aprecia que, a medida que el tamaño de la muestra crece, el poder estadístico también crece asintóticamente a 1, valor que equivale a tener la certeza de rechazar la hipótesis nula si esta es falsa.

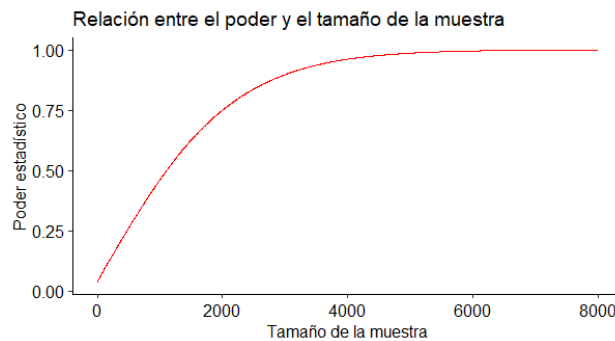


Figura 6.4: aumento del poder estadístico a medida que crece el tamaño de la muestra (manteniendo fijos el tamaño del efecto y el nivel de significación).

Script 6.2: aumento del poder estadístico a medida que crece el tamaño de la muestra.

```
1 library(ggpubr)
2
3 # Generar un vector con un rango para el tamaño de la muestra.
4 n <- seq(5, 8000, 5)
```



```

5
6 # Definir constantes
7 desv_est <- 6
8 alfa <- 0.05
9 tam_efecto <- 0.5
10
11 # Se calcula el poder con que se detecta el tamaño del efecto para
12 # cada tamaño de la muestra, asumiendo una prueba bilateral para
13 # una sola muestra.
14 poder <- power.t.test(n = n,
15                       delta = tam_efecto,
16                       sd = desv_est,
17                       sig.level = alfa,
18                       type = "two.sample",
19                       alternative = "two.sided")$power
20
21 # Crear un data frame.
22 datos <- data.frame(n, poder)
23
24 # Graficar la curva de poder.
25 g <- ggplot(datos, aes(n, poder))
26 g <- g + geom_line(colour = "red")
27 g <- g + ylab("Poder estadístico")
28 g <- g + xlab("Tamaño de la muestra")
29 g <- g + theme_pubr()
30 g <- g + ggtitle("Relación entre el poder y el tamaño de la muestra")
31
32 print(g)

```

6.4 CÁLCULO TEÓRICO DEL PODER

Como ya hemos mencionado a lo largo de este capítulo, el poder es la probabilidad de correctamente rechazar H_0 cuando es falsa, lo que equivale a la probabilidad de distinguir un efecto real de una mera casualidad. Ahora veremos algunos ejemplos de cómo podemos usar el poder.

Lola Drones, estudiante de computación, ha diseñado dos nuevos algoritmos (A y B) que resuelven un mismo problema como parte de su trabajo de titulación. Lola desea saber si existe diferencia entre los tiempos de ejecución de ambos algoritmos. Para ello, ha decidido realizar una prueba t con muestras pareadas, con un nivel de significación $\alpha = 0,05$, usando para ello 36 instancias del problema de tamaño fijo que se ejecutan bajo iguales condiciones con cada algoritmo. Además, Lola ya sabe que la diferencia en el tiempo de ejecución sigue una distribución normal con desviación estándar $\sigma = 12$ milisegundos. Así, Lola ha formulado las siguientes hipótesis:

H_0 : $\mu_{(A_i - B_i)} = 0$, es decir que la media de las diferencias en el tiempo de ejecución necesitado por los algoritmos A y B , para cada posible instancia i , es cero

H_A : $\mu_{(A_i - B_i)} \neq 0$

La figura 6.5 muestra cómo sería la distribución de la muestra (media de las diferencias en los tiempos de ejecución) si la hipótesis nula (H_0) fuese cierta, con las áreas correspondientes a la región de rechazo de H_0 coloreadas.

Supongamos por un momento que, en realidad, el algoritmo B es en promedio 4 milisegundos más rápido que el algoritmo A . En este caso tendríamos que la media de las diferencias es de -4 [ms], correspondiente

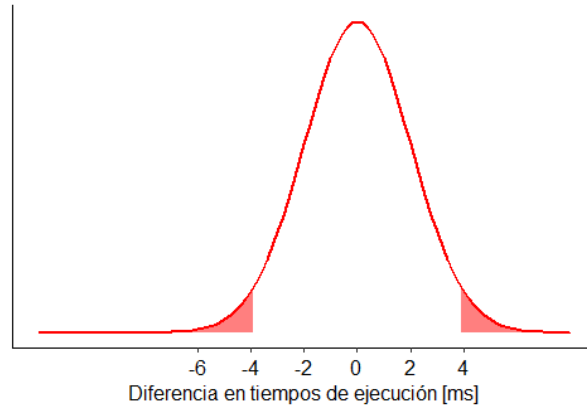


Figura 6.5: distribución de la diferencia de medias del tiempo de ejecución, señalando zonas de rechazo de la hipótesis nula.

al tamaño del efecto. En este caso, su distribución sería como muestra la figura 6.6 (ver script 6.3) en color azul. Al superponer esta nueva curva a la que ya teníamos bajo el supuesto de que la hipótesis nula fuera verdadera, vemos que el área de la curva real que se situaría dentro de la región de rechazo de la curva teórica es aquella coloreada en azul. Esta área corresponde al poder de la prueba t , que en este caso es de 0,516 de acuerdo al análisis teórico (ver script 6.1, líneas 77–86). Puesto que el poder corresponde a la probabilidad de **no** cometer un error de tipo II, de acuerdo al resultado obtenido se tiene que $\beta = 0,484$. ¡Lola no sería capaz de detectar una diferencia de -4 [ms] casi la mitad de las veces!

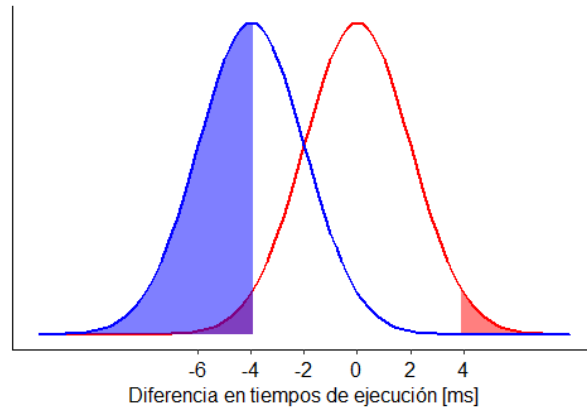


Figura 6.6: región de rechazo de la hipótesis nula en la distribución cuando el programa B es, en promedio, 4 milisegundos más rápido que el programa A .

Script 6.3: cálculo teórico del poder.

```
1 library(ggpubr)
2 library(pwr)
3
4 # Fijar valores conocidos.
5 sigma <- 12
6 alfa <- 0.05
7 n <- 36
8
9 # Calcular el error estándar.
10 SE <- sigma / sqrt(n)
11
12 # Gráficar la distribución muestral de la media de las diferencias si
```

```

13 # la hipótesis nula fuera verdadera.
14 x <- seq(-6 * SE, 4 * SE, 0.01)
15 y <- dnorm(x, mean = media_nula, sd = SE)
16 g <- ggplot(data = data.frame(x, y), aes(x))
17
18 g <- g + stat_function(
19   fun = dnorm,
20   args = list(mean = media_nula, sd = SE),
21   colour = "red", size = 1)
22
23 g <- g + ylab("")
24 g <- g + scale_y_continuous(breaks = NULL)
25 g <- g + scale_x_continuous(name = "Diferencia en tiempos de ejecución [ms]",
26                             breaks = seq(-6, 4, 2))
27
28 g <- g + theme_pubr()
29
30 # Colorear la región de rechazo de la hipótesis nula.
31 media_nula <- 0
32 Z_critico <- qnorm(alfa/2, mean = media_nula, sd = SE, lower.tail = FALSE)
33 q_critico_inferior <- media_nula - Z_critico
34 q_critico_superior <- media_nula + Z_critico
35
36 g <- g + geom_area(data = subset(df, x < q_critico_inferior),
37                   aes(y = y),
38                   colour = "red",
39                   fill = "red",
40                   alpha = 0.5)
41
42 g <- g + geom_area(data = subset(df, x > q_critico_superior),
43                   aes(y = y),
44                   colour = "red",
45                   fill = "red",
46                   alpha = 0.5)
47
48 print(g)
49
50 # Superponer la distribución muestral de la media de las diferencias
51 # si la la diferencia de medias fuera -4.
52 g <- g + stat_function(
53   fun = dnorm,
54   args = list(mean = media_efecto, sd = SE),
55   colour = "blue", size = 1)
56
57 # Colorear la región de la nueva curva situada en la región de
58 # rechazo de la curva original.
59 x1 <- seq(-6 * SE, 4 * SE, 0.01)
60 y1 <- dnorm(x, mean = media_efecto, sd = SE)
61 g <- g + geom_area(data = subset(data.frame(x1, y1),
62                                     x < q_critico_inferior),
63                   aes(x = x1, y = y1),
64                   colour = "blue",
65                   fill = "blue",
66                   alpha = 0.5)
67
68 g <- g + geom_area(data = subset(data.frame(x1, y1),
69                                     x > q_critico_superior),
70                   aes(x = x1, y = y1),
71                   colour = "blue",

```

```

72         fill = "blue",
73         alpha = 0.5)
74 print(g)
75
76 # Calcular el poder de acuerdo al análisis teórico.
77 poder <- pnorm(q_critico_inferior,
78               mean = media_efecto,
79               sd = SE,
80               lower.tail = TRUE)
81 + pnorm(q_critico_superior,
82         mean = media_efecto,
83         sd = SE,
84         lower.tail = FALSE)
85
86 cat("Poder = ", poder, "\n")
87
88 # Calcular la probabilidad de cometer un error tipo II.
89 beta <- 1 - poder_teorico
90 cat("Beta = ", beta, "\n")

```

6.5 CÁLCULO DEL PODER EN R

Desde luego, si trabajamos con R, podemos usar funciones para calcular el poder. Como primera alternativa, R trae incorporada la función `power.t.test(n, delta, sd, sig.level, power, type, alternative)` (empleada en los scripts 6.1 y 6.2), donde:

- **n**: tamaño de la muestra (por cada grupo, si corresponde).
- **delta**: diferencia observada entre las medias, o entre la media muestral y el valor nulo, no estandarizada.
- **sd**: desviación estándar observada.
- **sig.level**: nivel de significación.
- **power**: poder de la prueba.
- **type**: tipo de prueba t de Student (“two.sample” para diferencia de medias, “one.sample” para una sola muestra o “paired” para dos muestras pareadas).
- **alternative**: tipo de hipótesis alternativa (“one.sided” si es unilateral, “two.sided” si es bilateral).

Esta función entrega como resultado un objeto con diversos elementos (que podemos indexar del mismo modo que las columnas de una matriz de datos), entre los que se incluyen los 5 primeros argumentos definidos para la función.

Si revisamos con detenimiento los argumentos de la función `power.t.test()`, veremos que recibe el poder como uno de sus argumentos! Esto no parece tener sentido... ¿o sí?

Como ya hemos visto existe una relación entre: poder, tamaño de la muestra, tamaño del efecto y nivel de significación. A esta combinación de elementos debemos añadir también la desviación estándar, aunque no estudiaremos las matemáticas subyacentes.

En realidad, para usar la función `power.t.test()` siempre debemos señalar el tipo de prueba t con el que estamos trabajando y si la hipótesis alternativa es de una o dos colas. Esta función nos permite calcular cualquiera de los demás argumentos (tamaño de la muestra, tamaño del efecto, desviación estándar, nivel de significación o poder estadístico) para la prueba en cuestión a partir de los 4 argumentos restantes. Así, al argumento que queremos calcular se le asigna el valor `NULL` en la llamada.

Recordemos que en el ejemplo de la sección anterior, Lola Drones desea usar una prueba t bilateral para dos muestras pareadas a fin de determinar si hay diferencia entre los tiempos de ejecución promedio de ambos

algoritmos. Para ello, ha considerado $n = 36$ y $\alpha = 0,05$, sabiendo que $sd = 12$ [ms]. Las líneas 4 a 14 del script 6.4 muestran cómo calcular el poder para este ejemplo si se desea detectar un tamaño del efecto (δ) de 4 [ms], obteniéndose como resultado que el poder es de 0.494 (y $\beta = 1 - \text{poder} = 0,506$), ligeramente diferente al obtenido en forma teórica debido a errores de redondeo.

¿Cuántas instancias debería usar Lola para lograr un poder de 0,9, manteniendo $\alpha = 0,05$, $sd = 12$ [ms] y $\delta = 4$ [ms]? Las líneas 17–28 del script 6.4 muestran cómo hacer este cálculo, obteniéndose como resultado $n = 97$. Como el tamaño de la muestra siempre debe ser un entero positivo, la línea 27 aproxima el resultado al entero superior.

Otra alternativa es usar la función `pwr.t.test(n, d, sig.level, power, type, alternative)` (ver script 6.4, líneas 37–63), incluida en el paquete `pwr`, donde:

- **n**: tamaño de la muestra (por cada grupo, si corresponde).
- **d**: tamaño del efecto (d de Cohen).
- **sig.level**: nivel de significación.
- **power**: poder de la prueba.
- **type**: tipo de prueba t de Student (“two.sample” para diferencia de medias, “one.sample” para una sola muestra o “paired” para dos muestras pareadas).
- **alternative**: tipo de hipótesis alternativa (“greater” o “less” si es unilateral, “two.sided” si es bilateral).

Debemos fijarnos en que, si bien esta función opera de manera similar a `power.t.test()`, en este caso la desviación estándar y la diferencia son reemplazadas por el tamaño del efecto que podemos cuantificar, como ya vimos, mediante la d de Cohen. Sin embargo, debemos tener cuidado, pues la función `pwr.t.test()` solo es adecuada para una muestra, dos muestras pareadas o cuando ambas muestras tienen igual tamaño. En el caso de la prueba t para dos muestras independientes con diferentes tamaños, debemos usar, en cambio, la función `pwr.t2n.test(n1, n2, d, sig.level, power, alternative)`.

Script 6.4: cálculo del poder en R.

```

1 library(pwr)
2
3 # Fijar valores conocidos.
4 n <- 36
5 diferencia <- 4
6 desv_est <- 12
7 alfa <- 0.05
8 poder <- 0.9
9
10 # Calcular el poder usando la función power.t.test().
11 cat("Cálculo del poder con power.t.test()\n")
12
13 resultado <- power.t.test(n = n,
14                           delta = diferencia,
15                           sd = desv_est,
16                           sig.level = alfa,
17                           power = NULL,
18                           type = "paired",
19                           alternative = "two.sided")
20
21 print(resultado)
22
23 # Cálculo del tamaño de la muestra usando la función power.t.test().
24 cat("Cálculo del tamaño de la muestra con power.t.test()\n")
25
26 resultado <- power.t.test(n = NULL,
27                           delta = diferencia,
28                           sd = desv_est,
29                           sig.level = alfa,
```

```

30         power = poder,
31         type = "paired",
32         alternative = "two.sided")
33
34 n <- ceiling(resultado[["n"]])
35 cat("n = ", n, "\n")
36
37 # Calcular el tamaño del efecto (d de Cohen).
38 d <- (4 / desv_est) * ((n - 2) / (n - 1.25))
39
40 # Calcular el poder usando la función pwr.t.test().
41 cat("\n\nCálculo del poder con pwr.t.test()\n")
42
43 resultado <- pwr.t.test(n = n,
44                         d = d,
45                         sig.level = alfa,
46                         power = NULL,
47                         type = "paired",
48                         alternative = "two.sided")
49
50 print(resultado)
51
52 # Cálculo del tamaño de la muestra usando la función pwr.t.test().
53 cat("\n\nCálculo del tamaño de la muestra con pwr.t.test()\n")
54
55 resultado <- pwr.t.test(n = NULL,
56                         d = d,
57                         sig.level = alfa,
58                         power = poder,
59                         type = "paired",
60                         alternative = "two.sided")
61
62 n <- ceiling(resultado[["n"]])
63 cat("n = ", n, "\n")

```

6.6 EJERCICIOS PROPUESTOS

1. Define con tus propias palabras lo que es el tamaño del efecto.
2. Un estudio sobre el tiempo que necesitan los estudiantes para resolver una guía de ejercicios de Cálculo I, comparó un grupo de estudiantes que cursaban la asignatura por primera vez con un grupo que la cursaba en segunda ocasión. Sabiendo que este tiempo se distribuye normalmente en ambos casos, con varianzas similares, dibuja cómo se verían los datos si el efecto de repetir la asignatura sobre el tiempo requerido para resolver la guía fuera “grande” y si este efecto fuera “pequeño, pero positivo”.
3. Investiga cómo se calcula y cómo se interpreta la medida g de Hedges para el tamaño del efecto, e indica en qué casos es adecuada.
4. ¿Por qué se necesita conocer el tamaño del efecto?
5. ¿Cómo se relaciona el tamaño del efecto con la significación estadística?
6. ¿Por qué sería útil determinar un tamaño de muestra apropiado?
7. Explica en tus palabras lo que se muestra en la figura 6.4.
8. Ante algunas acusaciones de colusión, el Tribunal de la Libre Competencia quiere estudiar dos compañías del mercado de los seguros de automóviles. En base a datos del gremio de las aseguradoras, se puede asumir que el precio de las primas estándares para diferentes marcas de vehículos sigue una distribución

aproximadamente normal con desviación estándar de \$16.000. Fija los otros parámetros del estudio y determina qué tamaño debería tener la muestra de automóviles para detectar una diferencia de \$10.000 en el precio medio de las compañías bajo sospecha.

CAPÍTULO 7. INFERENCIA CON PROPORCIONES MUESTRALES

En el capítulo 5 conocimos las pruebas Z y t de Student para contrastar hipótesis con una y dos medias. Ahora estudiaremos los métodos de Wald y de Wilson para inferir acerca de una y dos proporciones, basándonos para ello en los textos de Diez y col. (2017, pp. 274-286), NIST/SEMATECH (2013, pp. 7.2.4, 7.2.4.1), Pértiga y Pita (2004), Champely, Ekstrom, Dalgaard, Gill, Weibelzahl, Anandkumar, Ford, Volcic y de Rosario (2020) y Kabacoff (2017).

7.1 MÉTODO DE WALD

En el capítulo 3 vimos que, cuando queremos responder preguntas del tipo “¿qué proporción de la ciudadanía apoya al gobierno actual?”, estamos hablando de una variable aleatoria que sigue una distribución binomial. En general, no conocemos la **probabilidad de éxito** p de la población, por lo que tenemos que usar el estimador puntual (correspondiente a la proporción de éxito de la muestra), denotado por \hat{p} . Este estimador se distribuye de manera cercana a la normal cuando se cumplen las siguientes condiciones:

1. Las observaciones de la muestra son independientes.
2. Se cumple la **condición de éxito-fracaso**, que establece que se espera observar al menos 10 observaciones correspondientes a éxito y al menos 10, correspondientes a fracasos. Matemáticamente, $np \geq 10$ y $n(1 - p) \geq 10$.

Así, si la distribución muestral de \hat{p} cumple con las condiciones anteriores, se dice que es cercana a la normalidad con media $\mu = p$, desviación estándar $\sigma = \sqrt{p(1 - p)}$ y error estándar dado por la ecuación 7.1.

$$SE_{\hat{p}} = \sqrt{\frac{p(1 - p)}{n}} \quad (7.1)$$

7.1.1 Método de Wald para una proporción

El **método de Wald** permite construir intervalos de confianza y contrastar hipótesis bajo el supuesto de normalidad para una proporción. Consideremos el siguiente ejemplo: Aquiles Baeza, ingeniero en informática, desea conocer qué proporción de las ejecuciones de un algoritmo de ordenamiento para instancias con 100.000 elementos (bajo iguales condiciones de hardware y sistema) tardan menos de 25 segundos. Para ello, registró los tiempos de ejecución para 150 instancias generadas de manera aleatoria, encontrando que 64 % de dichas instancias fueron resueltas en un tiempo menor al señalado.

Si bien no conocemos la probabilidad real de éxito para la población, sabemos que $\hat{p} = 0,64$. Así, si se cumplen las condiciones para que la distribución de \hat{p} sea cercana a la normal, podemos construir un intervalo de confianza para la verdadera proporción muestral.

En el enunciado del ejemplo nos indican que las instancias del problema fueron escogidas de manera aleatoria y sabemos que éstas representan menos del 10 % del total de instancias posibles, con lo que se verifica la

independencia de las observaciones. Por otra parte, nos dicen que la proporción de éxito es $\hat{p} = 0,64$, por lo que esperamos encontrar $0,64 \cdot 150 = 96$ instancias que tardan menos de 25 segundos y $(1 - 0,64) \cdot 150 = 54$ fracasos (instancias que tardan 25 segundos o más), con lo que se cumple la condición de éxito-fracaso. En consecuencia, podemos asumir que la distribución muestral de \hat{p} sigue aproximadamente a la normal.

Podemos estimar el error estándar usando la ecuación 7.1, reemplazando p por el estadístico \hat{p} :

$$SE = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} = \sqrt{\frac{0,64(1 - 0,64)}{150}} = 0,0392$$

Con ello, construimos el intervalo de confianza para un nivel de significación $\alpha = 0,05$ usando la ecuación general (4.6) con \hat{p} como estimador puntual:

$$\hat{p} \pm z^* \cdot SE \rightarrow 0,64 \pm 1,96 \cdot 0,0392 \rightarrow [0,5632; 0,7168]$$

Este intervalo significa que tenemos 95 % de confianza que la proporción de instancias (de 100.000 elementos) del problema que el algoritmo ordena en menos de 25 segundos se encuentra entre 56,32 % y 71,6 %.

Desde luego, también podemos usar el modelo normal en el contexto de la prueba de hipótesis para una proporción. Para ello, se deben cumplir las condiciones de independencia y éxito-fracaso que ya verificamos para construir el intervalo de confianza, pero en este caso tenemos que verificar la segunda condición con el valor nulo, denotado p_0 . Una vez verificadas ambas condiciones, el error estándar y el estadístico Z que permiten determinar el p-valor se calculan usando las ecuaciones 7.2 y 7.3, respectivamente.

$$SE = \sqrt{\frac{p_0(1 - p_0)}{n}} \quad (7.2)$$

$$Z = \frac{\hat{p} - p_0}{SE} \quad (7.3)$$

Supongamos ahora, volviendo a nuestro ejemplo, que Baeza afirma que más del 70 % de las instancias de tamaño 100.000 se ejecutan en menos de 25 segundos. Sin embargo, su jefe no está seguro, por lo que decide comprobarlo mediante una prueba de hipótesis con un nivel de significación $\alpha = 0,05$ (recordemos que $n = 150$ y $\hat{p} = 0,64$):

H_0 : el 70 % de las instancias se ejecutan en menos de 25 segundos.

H_A : más del 70 % de las instancias se ejecutan en menos de 25 segundos.

De acuerdo a las hipótesis formuladas por el jefe de Baeza, el valor nulo es $p_0 = 0,7$, con lo que estas pueden formularse matemáticamente como:

Denotando como p a la proporción de todas las instancias de tamaño 100.000 que se ejecutan en menos de 25 segundos y considerando el valor hipotético $p_0 = 0,7$ para este parámetro:

H_0 : $p = p_0$

H_A : $p > p_0$

Ya antes habíamos comprobado que se verifica la independencia de las observaciones. Además, considerando que el valor nulo fuese verdadero esperaríamos encontrar $0,7 \cdot 150 = 105$ éxitos y $(1 - 0,7) \cdot 150 = 45$ fracasos, ambos valores mayores que 10, por lo que la condición de éxito-fracaso se verifica.

Con ello, podemos calcular el estadístico de prueba:

$$SE = \sqrt{\frac{p_0(1 - p_0)}{n}} = \sqrt{\frac{0,7(1 - 0,7)}{150}} = 0,0374$$

$$Z = \frac{\hat{p} - p_o}{SE} = \frac{0,64 - 0,7}{0,0374} = -1,6043$$

El valor p asociado, calculado en R mediante la llamada a la función `pnorm(-1.6042, lower.tail = FALSE)`, es $p = 0,9456$. En consecuencia, la evidencia no es suficiente para rechazar la hipótesis nula, por lo que se concluye, con 95% de confianza, que no es cierto que el algoritmo se ejecute en menos de 25 segundos para más del 70% de las instancias de tamaño 100.000.

R no ofrece esta prueba, como función. Sin embargo, podemos hacerla como muestra el script 7.1 para nuestro ejemplo.

Script 7.1: método de Wald para una proporción.

```

1 # Fijar valores conocidos
2 n <- 150
3 p_exito <- 0.64
4 alfa <- 0.05
5 valor_nulo <- 0.7
6
7 # Construcción del intervalo de confianza.
8 error_est <- sqrt((p_exito * (1 - p_exito)) / n)
9 Z_critico <- qnorm(alfa / 2, lower.tail = FALSE)
10 inferior <- p_exito - Z_critico * error_est
11 superior <- p_exito + Z_critico * error_est
12 cat("Intervalo de confianza = [", inferior, ", ", superior, "]\n", sep = "")
13
14 # Prueba de hipótesis.
15 error_est_hip <- sqrt((valor_nulo * (1 - valor_nulo)) / n)
16 Z <- (p_exito - valor_nulo) / error_est_hip
17 p <- pnorm(Z, lower.tail = FALSE)
18 cat("Hipótesis alternativa unilateral\n")
19 cat("Z =", Z, "\n")
20 cat("p =", p)

```

7.1.2 Método de Wald para dos proporciones

También podemos usar el método de Wald para estudiar la **diferencia entre las proporciones** de dos poblaciones, considerando para ello como estimador puntual la diferencia $\hat{p}_1 - \hat{p}_2$.

De manera similar a lo que ya vimos para una única proporción, también en este caso debemos verificar ciertas condiciones antes de poder aplicar el modelo normal:

1. Cada proporción, por separado, sigue el modelo normal.
2. Las dos muestras son independientes una de la otra.

El error estándar para la diferencia entre dos proporciones muestrales está dado por la ecuación 7.4, donde p_1 y p_2 corresponden a las proporciones de las poblaciones, y n_1 y n_2 , a los tamaños de las muestras. La construcción del intervalo de confianza se realiza, una vez más, con la ecuación general 4.6.

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{SE_{\hat{p}_1}^2 + SE_{\hat{p}_2}^2} = \sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}} \quad (7.4)$$

A modo de ejemplo, supongamos que la Facultad de Ingeniería de una prestigiosa universidad desea determinar si la tasa de reprobación de estudiantes que rinden la asignatura de programación por primera vez es igual para hombres y mujeres. Para ello, se examina la situación final de los estudiantes que rindieron la asignatura durante el segundo semestre de 2017. Para una muestra de 48 hombres (de un total de 632), se encontró que 26 de ellos reprobaron la asignatura. De manera similar, para una muestra de 42 mujeres (de un total de 507), se encontraron 20 reprobaciones¹, con ambas muestras tomadas de manera aleatoria.

Como ya es habitual, comencemos por verificar las condiciones de normalidad para cada una de las muestras. En ambos casos, las observaciones son independientes entre sí, pues provienen de personas diferentes que representan a menos del 10 % de la población. Además, los datos entregados evidencian que en ambos casos se cumple la condición de éxito-fracaso. Adicionalmente, ambas muestras son independientes entre sí, pues ambas categorías se excluyen mutuamente. Con esto último se verifican entonces las condiciones de normalidad para la diferencia de proporciones.

Sean \hat{p}_1 y \hat{p}_2 las proporciones de éxito muestrales (considerando en este contexto la reprobación como éxito) para hombres y mujeres, respectivamente:

$$\hat{p}_1 = 26/48 = 0,5417$$

$$\hat{p}_2 = 20/42 = 0,4762$$

$$\hat{p}_1 - \hat{p}_2 = 0,5417 - 0,4762 = 0,0655$$

El error estándar puede estimarse como:

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} = \sqrt{\frac{0,5417(1-0,5417)}{48} + \frac{0,4762(1-0,4762)}{42}} = 0,1054$$

Suponiendo un nivel de significación $\alpha = 0,05$, el intervalo de confianza corresponde a:

$$\hat{p}_1 - \hat{p}_2 \pm z^* SE_{\hat{p}_1 - \hat{p}_2} \rightarrow 0,0655 \pm 1,96 \cdot 0,1054 \rightarrow [-0,1411; 0,2721]$$

En consecuencia, podemos afirmar con 95 % de confianza que la diferencia en la tasa de reprobación de la asignatura de programación para hombres y mujeres varía entre -14,11 % y 27,21 %.

Desde luego, también podemos realizar pruebas de hipótesis en este escenario. Para el ejemplo tenemos que:

H_0 : no hay diferencia en la tasa de reprobación de hombres y mujeres.

H_A : las tasas de reprobación son diferentes para hombres y mujeres.

Matemáticamente:

Denotando como p_1 y p_2 a las proporciones de hombres y mujeres, respectivamente, que reprobaban la asignatura de programación la primera vez que la cursan:

H_0 : $p_1 - p_2 = 0$

H_A : $p_1 - p_2 \neq 0$

Ya verificamos las condiciones para operar bajo el supuesto de normalidad cuando construimos el intervalo de confianza. Sin embargo, **cundo la hipótesis nula supone que no hay diferencia entre las proporciones**, la verificación de la condición de éxito-fracaso y la estimación del error estándar se realizan usando para ello la **proporción agrupada**, dada por la ecuación 7.5, donde $\hat{p}_1 n_1$ y $\hat{p}_2 n_2$ representan la cantidad de éxitos en la primera y segunda muestra, respectivamente.

¹Los datos aquí presentados son ficticios, creados únicamente con fines pedagógicos.

$$\hat{p} = \frac{\text{número de éxitos}}{\text{número de casos}} = \frac{\hat{p}_1 n_1 + \hat{p}_2 n_2}{n_1 + n_2} \quad (7.5)$$

Así, en este caso tenemos:

$$\hat{p} = \frac{\hat{p}_1 n_1 + \hat{p}_2 n_2}{n_1 + n_2} = \frac{0,5417 \cdot 48 + 0,4762 \cdot 42}{48 + 42} = 0,5111$$

En consecuencia, en el caso de los hombres esperamos encontrar $\hat{p}n_1 > 24$ éxitos (reprobaciones) y $(1 - \hat{p})n_1 > 23$ fracasos. Del mismo modo, para las mujeres esperamos $\hat{p}n_2 > 21$ éxitos y $(1 - \hat{p})n_2 > 20$ fracasos, con lo que se verifican las condiciones para emplear el modelo normal.

El error estándar se calcula, como ya mencionamos, usando la proporción agrupada:

$$SE = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n_1} + \frac{\hat{p}(1 - \hat{p})}{n_2}} = \sqrt{\frac{0,5111 \cdot (1 - 0,5111)}{48} + \frac{0,5111 \cdot (1 - 0,5111)}{42}} = 0,1056$$

El estimador puntual para la diferencia es $\hat{p}_1 - \hat{p}_2 = 0,0655$, con lo cual el estadístico de prueba está dado por:

$$Z = \frac{\text{estimador puntual} - \text{valor nulo}}{SE} = \frac{0,0655 - 0}{0,1056} = 0,6203$$

En consecuencia, el valor p correspondiente es $p = 0,5351$. Puesto que el valor p es mayor que $\alpha = 0,05$, se falla en rechazar la hipótesis nula. Así, podemos decir con 95 % de confianza que no existe evidencia suficiente para concluir que hay diferencia en la tasa de reprobación de hombres y mujeres para el primer curso de programación.

El script 7.2 muestra el desarrollo de este ejemplo en R.

Script 7.2: método de Wald para la diferencia entre dos proporciones (ejemplo 1).

```

1 # Fijar valores conocidos
2 n_hombres <- 48
3 n_mujeres <- 42
4 exitos_hombres <- 26
5 exitos_mujeres <- 20
6 alfa <- 0.05
7 valor_nulo <- 0
8
9 # Calcular probabilidades de éxito.
10 p_hombres <- exitos_hombres / n_hombres
11 p_mujeres <- exitos_mujeres / n_mujeres
12
13 # Estimar la diferencia.
14 diferencia <- p_hombres - p_mujeres
15
16 # Construcción del intervalo de confianza.
17 error_hombres <- (p_hombres * (1 - p_hombres)) / n_hombres
18 error_mujeres <- (p_mujeres * (1 - p_mujeres)) / n_mujeres
19 error_est <- sqrt(error_hombres + error_mujeres)
20 Z_critico <- qnorm(alfa / 2, lower.tail = FALSE)
21 inferior <- diferencia - Z_critico * error_est
22 superior <- diferencia + Z_critico * error_est
23 cat("Intervalo de confianza = [", inferior, ", ", superior, "]\n", sep = "")
24
25 # Prueba de hipótesis.

```

```

26 p_agrupada <- ( exitos_hombres + exitos_mujeres ) / ( n_hombres + n_mujeres )
27 error_hombres <- ( p_agrupada * ( 1 - p_agrupada ) ) / n_hombres
28 error_mujeres <- ( p_agrupada * ( 1 - p_agrupada ) ) / n_mujeres
29 error_est_hip <- sqrt( error_hombres + error_mujeres )
30 Z <- ( diferencia - valor_nulo ) / error_est_hip
31 p <- 2 * pnorm( Z, lower.tail = FALSE )
32 cat( "Hipótesis alternativa bilateral\n" )
33 cat( "Z =", Z, "\n" )
34 cat( "p =", p )

```

Cuando contrastamos hipótesis para la **diferencia entre dos proporciones con un valor nulo distinto de 0**, el procedimiento es ligeramente diferente. En este caso, la comprobación de la condición de éxito-fracaso se realiza de manera independiente para ambas muestras y el error estándar se calcula, como ya se estudió para los intervalos de confianza, mediante la ecuación 7.4.

Supongamos ahora que la Facultad de Ingeniería de la Universidad anterior ha decidido replicar el estudio realizado para el curso de programación, esta vez para una asignatura de física. No obstante, las autoridades están convencidas de que la tasa de reprobación es 10 % mayor para los hombres y que, incluso, la diferencia podría ser mayor. Desean comprobar con un nivel de confianza de 95 % y para ello, seleccionaron aleatoriamente a 89 de los 1.023 hombres y a 61 de las 620 mujeres de la cohorte correspondiente al primer semestre de 2019. En las muestras se encuentran, respectivamente, 45 y 21 reprobaciones.

Las hipótesis son, en este caso:

H_0 : la tasa de reprobación de los hombres es exactamente 10 % más alta que la de las mujeres.

H_A : la tasa de reprobación de los hombres es más de 10 % más alta que la de las mujeres.

Matemáticamente:

Denotando como p_1 y p_2 a las proporciones de hombres y mujeres, respectivamente, que reprueban la asignatura de física estudiada la primera vez que la cursan:

H_0 : $p_1 - p_2 = 0,1$

H_A : $p_1 - p_2 > 0,1$

Al igual que en los ejemplos previos, las observaciones de cada muestra son independientes entre sí pues corresponden a menos del 10 % de la población y fueron escogidos aleatoriamente. A su vez, los datos proporcionados indican que se cumple la condición de éxito-fracaso para cada muestra. Como ambas muestras pertenecen a grupos diferentes de estudiantes, son independientes entre sí. En consecuencia, se cumplen las condiciones para operar bajo el modelo normal.

En el caso de los hombres, la tasa de éxito se estima como:

$$\hat{p}_1 = \frac{45}{89} = 0,5056$$

Análogamente, para las mujeres tenemos:

$$\hat{p}_2 = \frac{21}{61} = 0,3443$$

Con lo que el estimador puntual para la diferencia es:

$$\hat{p}_1 - \hat{p}_2 = 0,5056 - 0,3443 = 0,1613$$

Ahora calculamos el error estándar:

$$SE = \sqrt{\frac{0,5056(1 - 0,5056)}{89} + \frac{0,3443(1 - 0,3443)}{61}} = 0,0807$$

Con lo cual podemos calcular el estadístico de prueba:

$$Z = \frac{\text{estimador puntual} - \text{valor nulo}}{SE} = \frac{0,1613 - 0,1}{0,0807} = 0,7596$$

Con lo que se puede obtener el valor p , correspondiente a $p = 0.2237 > \alpha = 0,05$.

En consecuencia, se falla en rechazar H_0 en favor de H_A , por lo que concluimos, con 95% de confianza, que la tasa de reprobación de los hombres es 10% superior a la de las mujeres para el curso de física.

En R, esta prueba puede realizarse como muestra el script 7.3.

Script 7.3: método de Wald para la diferencia entre dos proporciones (ejemplo 2).

```
1 # Fijar valores conocidos
2 n_hombres <- 89
3 n_mujeres <- 61
4 exitos_hombres <- 45
5 exitos_mujeres <- 21
6 alfa <- 0.05
7 valor_nulo <- 0.1
8
9 # Calcular probabilidades de éxito.
10 p_hombres <- exitos_hombres / n_hombres
11 p_mujeres <- exitos_mujeres / n_mujeres
12
13 # Estimar la diferencia.
14 diferencia <- p_hombres - p_mujeres
15
16 # Prueba de hipótesis.
17 p_agrupada <- (exitos_hombres + exitos_mujeres) / (n_hombres + n_mujeres)
18 error_hombres <- (p_hombres * (1 - p_hombres)) / n_hombres
19 error_mujeres <- (p_mujeres * (1 - p_mujeres)) / n_mujeres
20 error_est <- sqrt(error_hombres + error_mujeres)
21 Z <- (diferencia - valor_nulo) / error_est
22 p <- pnorm(Z, lower.tail = FALSE)
23 cat("Hipótesis alternativa bilateral\n")
24 cat("Z =", Z, "\n")
25 cat("p =", p)
```

7.2 MÉTODO DE WILSON

El método de Wald, tratado en la sección anterior, es el método que tradicionalmente se ha usado y el que aparece en la mayoría de los libros clásicos de inferencia estadística. Sin embargo, el método está siendo muy criticado hoy en día debido a que hace importantes simplificaciones matemáticas en su procedimiento y ya hay evidencia empírica que ha demostrado sus limitaciones (Agresti & Coull, 1998).

Gracias al aumento del poder de cómputo y la disponibilidad de software estadístico, han surgido diversas alternativas, entre las cuales destaca el **método de Wilson** (junto con algunas variaciones), considerado el más robusto por diversos autores (Agresti & Coull, 1998; Brown y col., 2001; Devore, 2008; Wallis, 2013). Este método opera del mismo modo que el de Wald, aunque las fórmulas empleadas para estimar la proporción en la muestra y el error estándar son diferentes.

En R, podemos hacer esta prueba usando la función `prop.test(x, n, p, alternative, conf.level, ...)`, cuyos principales parámetros son:

- `x`: cantidad de éxitos en la muestra.
- `n`: tamaño de la muestra.
- `p`: valor nulo (por defecto, `p=NULL`).
- `alternative`: tipo de hipótesis alternativa, por defecto bilateral (`alternative="two.sided"`), y valores `"less"` y `"greater"` para hipótesis unilaterales.
- `conf.level`: nivel de confianza (`conf.level=0.95` por defecto).

El script 7.4 muestra el uso de esta función con el mismo ejemplo que usamos para presentar la prueba de Wald para una proporción. Del mismo modo, el script 7.5 usa la función `prop.test()` para el primer ejemplo del método de Wald para la diferencia entre dos proporciones. Sin embargo, esta función tiene la limitante de que, al trabajar con dos proporciones, no permite establecer un valor nulo distinto de cero para la diferencia.

Script 7.4: método de Wilson para una proporción.

```
1 # Fijar valores conocidos
2 n <- 150
3 p_exito <- 0.64
4 alfa <- 0.05
5 valor_nulo <- 0.7
6
7 # Calcular cantidad de éxitos.
8 exitos <- p_exito * n
9
10 # Prueba de Wilson en R.
11 prueba <- prop.test(exitos, n = n, p = valor_nulo,
12                     alternative = "greater", conf.level = 1 - alfa)
13
14 print(prueba)
```

Script 7.5: método de Wilson para la diferencia entre dos proporciones.

```
1 # Fijar valores conocidos (hombres, mujeres)
2 n <- c(48, 42)
3 exitos <- c(26, 20)
4 alfa <- 0.05
5
6 # Prueba de Wilson en R.
7 prueba <- prop.test(exitos, n = n, alternative = "two.sided",
8                     conf.level = 1 - alfa)
9
10 print(prueba)
```

7.3 PODER Y PRUEBAS DE PROPORCIONES

En el capítulo anterior conocimos el poder estadístico y vimos que está relacionado con el nivel de significación, el tamaño de la muestra y el tamaño del efecto que queremos detectar.

R base nos ofrece la función `power.prop.test(n, p1, p2, sig.level, power, alternative)`, donde:

- `n`: número de observaciones por cada grupo.
- `p1`: probabilidad de éxito en un grupo.
- `p2`: probabilidad de éxito en otro grupo.

- **sig.level**: nivel de significación.
- **power**: poder de la prueba.
- **alternative**: tipo de hipótesis alternativa (“one.sided” si es unilateral, “two.sided” si es bilateral).

Al igual que vimos en el capítulo anterior para la función `power.t.test()`, recibe cuatro de los primeros argumentos y al restante debe asignársele el valor `NULL`. Como resultado, retorna un objeto que incluye el valor calculado para el argumento faltante.

Una vez más, el paquete `pwr` de R nos ofrece varias funciones que podemos usar como alternativa:

- `pwr.p.test(h, n, sig.level, power, alternative)`: para pruebas con una única proporción.
- `pwr.2p.test(h, n, sig.level, power, alternative)`: para pruebas con dos proporciones donde ambas muestras son de igual tamaño.
- `pwr.2p2n.test(h, n1, n2, sig.level, power, alternative)`: para pruebas con dos proporciones y muestras de diferente tamaño.

Donde:

- **h**: tamaño de efecto.
- **n, n1, n2**: tamaño(s) de la(s) muestra(s).
- **sig.level**: nivel de significación.
- **power**: poder.
- **alternative**: tipo de hipótesis alternativa (“two.sided”, “less” o “greater”).

El funcionamiento de esta familia de funciones es igual al que ya conocimos en el capítulo anterior para la función `pwr.t.test()`. Se entrega el parámetro **alternative** y todos los demás excepto uno (al cual debe asignarse explícitamente el valor `NULL`). Como resultado, la función calcula dicho valor.

El tamaño del efecto puede calcularse como muestra la ecuación 7.6, implementada en R en la función `ES.h(p1, p2)` del paquete `pwr`.

$$h = 2 \arcsin(\sqrt{p_1}) - 2 \arcsin(\sqrt{p_2}) \quad (7.6)$$

En el caso de una única proporción, los autores del paquete `pwr` sugieren usar $p_2 = 0,5$ (Champely y col., 2020).

Otra función que nos puede ser de ayuda es `bsamsize(p1, p2, fraction, alpha, power)`, del paquete `Hmisc`. En el caso de una prueba de Wilson con dos muestras, calcula los tamaños de cada grupo dados los siguientes argumentos:

- **p1**: probabilidad de la población para el grupo 1.
- **p2**: probabilidad del grupo 2.
- **fraction**: fracción de las observaciones en el grupo 1 ($n1/(n1 + n2)$).
- **alpha**: nivel de significación.
- **power**: poder deseado.

7.4 EJERCICIOS PROPUESTOS

1. ¿En qué condiciones la distribución muestral de una proporción tiene comportamiento aproximadamente normal?
2. ¿Cómo se calcula la desviación estándar de la distribución muestral de las proporciones bajo estas condiciones (según el método de Wald)?
3. ¿Cómo se calcula un intervalo de confianza para la verdadera proporción (según el método de Wald)?

4. El patrón de un gran fundo de nogales está preocupado porque se ha detectado la presencia de una plaga en varios árboles. Si bien existe un pesticida para el parásito, este es bastante caro y su aplicación solo se justifica económicamente si más del 20 % de los árboles está infectado. En consecuencia, el patrón ha decidido estimar la extensión de la infestación revisando una muestra aleatoria de 200 nogales (una porción bastante pequeña de los más de 20.000 árboles en el fundo). En base a lo anterior, determina:
 - a) ¿Cuál es la variable dicotómica (experimento Bernulli) en este caso?
 - b) ¿Cuál es el parámetro de interés?
 - c) ¿Qué estimador existe para este parámetro?
 - d) ¿Qué hipótesis respondería las dudas del patrón del fundo?
5. En el experimento del ejercicio anterior se encontró que 45 árboles de la muestra estaban infectados:
 - a) ¿Se puede asumir que esta proporción muestral sigue el modelo normal?
 - b) Independientemente de la respuesta anterior, obtén un intervalo con 95 % confianza para la verdadera proporción de árboles infectados en el fundo.
 - c) ¿Qué recomendarías al patrón del fundo?
6. Como el patrón sigue con dudas, ahora pregunta: ¿cuántos árboles debería revisar en una muestra para estar 99 % confiado que más del 20 % de los árboles están infectados, con solo 10 % de probabilidades de equivocarse si la verdadera proporción fuera 18 %? ¿Cómo se puede calcular esto? ¿Cuál debiera ser la respuesta a la pregunta del patrón?
7. ¿En qué condiciones la distribución muestral de la diferencia de dos proporciones tiene comportamiento aproximadamente normal?
8. ¿Cómo se calcula el error estándar de la diferencia entre dos proporciones (según el método de Wald)?
9. ¿Cómo se calcula un intervalo de confianza para la verdadera diferencia entre dos proporciones (según el método de Wald)?
10. Un laboratorio homeopático acaba de lanzar un tónico que asegura que ayuda a prevenir el resfrío durante el periodo invernal, con igual eficacia tanto en mujeres como en hombres. Para comprobar esta promesa, el laboratorio está realizando un estudio de la eficacia del producto en una muestra aleatoria de 100 mujeres y 200 hombres:
 - a) ¿Cuál es el parámetro de interés y qué estimador se podría usar?
 - b) ¿Qué hipótesis se deberían docimar para comprobar o refutar la homogeneidad de la eficacia del tónico para el resfrío?
11. El estudio anterior encontró que, durante las semanas de prueba, 38 mujeres y 102 hombres presentaron síntomas de resfrío. ¿Es homogénea la eficacia del producto con un nivel de significación de 0,05?
12. ¿Qué poder tuvo la prueba anterior?
13. ¿Qué tamaño deberían tener las muestras aleatorias de mujeres y hombres (manteniendo la proporción del ejemplo) para conseguir un poder de 0,85 con 99 % de confianza?
14. Las fórmulas presentadas en la sección 7.1, se conocen colectivamente como “Método de Wald”, el que ya no es recomendado por académicos del área. Usando la bibliografía citada, ¿cuáles son las fórmulas del método de Wilson para estimar el error estándar de la proporción y su extensión a la prueba de hipótesis e intervalos de confianza?
15. Investiga para qué sirve y cómo funciona el parámetro `correct` (verdadero por defecto) de la función `prop.test()` de R.

REFERENCIAS

- Agresti, A. & Coull, B. A. (1998).
Approximate is better than “exact” for interval estimation of binomial proportions.
The American Statistician, 52(2), 119-126.
- Brown, L. D., Cai, T. T. & DasGupta, A. (2001). Interval estimation for a binomial proportion.
Statistical science, 16(2), 101-117.
- Champely, S., Ekstrom, C., Dalgaard, P., Gill, J., Weibelzahl, S., Anandkumar, A., Ford, C., Volcic, R. & de Rosario, H. (2020). *pwr: Basic Functions for Power Analysis*. Consultado el 1 de octubre de 2021, desde <https://cran.r-project.org/web/packages/pwr/vignettes/pwr-vignette.html>
- Devore, J. L. (2008). *Probabilidad y Estadística para Ingeniería y Ciencias* (7.^a ed.). CENAGE Learning.
- Diez, D., Barr, C. D. & Çetinkaya-Rundel, M. (2017). *OpenIntro Statistics* (3.^a ed.).
<https://www.openintro.org/book/os/>.
- Freund, R. J. & Wilson, W. J. (2003). *Statistical Methods* (2.^a ed.). Academic Press.
- Kabacoff, R. I. (2017). *Power Analysis*.
Consultado el 1 de octubre de 2021, desde <https://www.statmethods.net/stats/power.html>
- Kassambara, A. (2019). *T-test Effect Size using Cohen’s d Measure*.
Consultado el 27 de abril de 2021, desde <https://www.datanovia.com/en/lessons/t-test-effect-size-using-cohens-d-measure/#cohens-d-for-paired-samples-t-test>
- NIST/SEMATECH. (2013). *e-Handbook of Statistical Methods*.
Consultado el 29 de abril de 2021, desde <http://www.itl.nist.gov/div898/handbook/>
- Pértega, S. & Pita, S. (2004).
Asociación de variables cualitativas: El test exacto de Fisher y el test de McNemar. Consultado el 29 de abril de 2021, desde <https://www.fisterra.com/mbe/investiga/fisher/fisher.asp#McNemar>
- Wallis, S. (2013). Binomial confidence intervals and contingency tests: mathematical fundamentals and the evaluation of alternative methods. *Journal of Quantitative Linguistics*, 20(3), 178-208.